

# MAPGD: Multi-Agent Prompt Gradient Descent for Collaborative Prompt Optimization

Anonymous ACL submission

## Abstract

Prompt engineering is crucial for fully leveraging large language models, yet existing prompt optimization methods often rely on a single refinement trajectory, leading to instability, conflicting update signals, and inefficient use of query budgets. We propose **Multi-Agent Prompt Gradient Descent (MAPGD)**, a gradient-inspired framework that coordinates multiple complementary prompt editing signals for robust and interpretable optimization. MAPGD decomposes prompt refinement into orthogonal dimensions (e.g., instruction clarity, example selection, format, and style) and aggregates textual pseudo-gradients via semantic embedding, conflict-aware clustering, and adaptive fusion. To improve robustness, we introduce HCGC, which enforces angular separation between conflicting directions, and CAAW, which dynamically calibrates editing contributions based on validation feedback. Experiments on diverse classification and reasoning benchmarks show that MAPGD achieves consistent gains in accuracy over strong baselines, while ablation results indicate that coordinated gradient fusion and adaptive weighting are critical to stable optimization. Together, these findings suggest that MAPGD offers a robust and interpretable framework for automatic prompt optimization.

## 1 Introduction

Large Language Models (LLMs), trained on vast web-scale corpora, have demonstrated remarkable generalization capabilities across a wide range of natural language processing (NLP) tasks, including classification, reasoning, and generation (Achiam et al., 2023; Bubeck et al., 2023). A critical factor underlying this performance is the design of prompts, which serve as the primary interface for steering model behavior. Despite their importance, prompts are still predominantly crafted through manual trial-and-error, requiring substantial hu-

man effort (Jiang et al., 2022) and domain expertise (Reynolds and McDonell, 2021; Zamfirescu-Pereira et al., 2023). This reliance on manual design limits scalability and motivates the development of automatic or semi-automatic prompt optimization methods that are both effective and interpretable.

Existing approaches to prompt optimization largely fall into two categories. The first line of work approximates gradient-based optimization through auxiliary models or differentiable prompt representations (Qin and Eisner, 2021; Deng et al., 2022). While theoretically appealing, these methods often assume access to internal model states or parameters (Shin et al., 2020; Lester et al., 2021), making them impractical in real-world, API-only settings. The second line of work formulates prompt optimization as a discrete search or reinforcement learning problem, using feedback-driven edits or Monte Carlo exploration (Zhang et al., 2022; Zhou et al., 2022). Although effective in certain cases, these methods can suffer from high computational cost, unstable optimization trajectories, and limited interpretability. Collectively, these limitations highlight the need for more robust, scalable, and transparent prompt optimization frameworks.

To bridge this gap, Pryzant et al. (Pryzant et al., 2023) proposed ProTeGi, a non-parametric framework that reinterprets prompt optimization as gradient descent in natural language space. ProTeGi generates textual pseudo-gradients from model errors and iteratively refines prompts using beam search and bandit-based selection, achieving strong empirical performance. However, ProTeGi relies on a single-agent optimization framework that aggregates heterogeneous refinement signals into a unified trajectory. Such entanglement may give rise to implicit objective conflicts, thereby impairing optimization stability and convergence quality.

A natural extension is to decompose prompt op-

084 timization into multiple specialized perspectives. 136  
085 Multi-agent formulations enable parallel explora- 137  
086 tion of heterogeneous refinement dimensions, 138  
087 increasing signal diversity and reducing prema- 139  
088 ture convergence. However, introducing multiple 140  
089 agents also makes semantic conflicts inevitable: dif- 141  
090 ferent agents may propose mutually incompatible 142  
091 updates, such as simultaneously expanding illustra- 143  
092 tive examples while enforcing concise instructions. 144  
093 Without explicit coordination mechanisms, such 145  
094 conflicts can destabilize optimization and under- 146  
095 mine the benefits of multi-agent collaboration. 147

096 In this work, we propose MAPGD, a unified 148  
097 framework that addresses these challenges by inte- 149  
098 grating multi-agent specialization with principled 150  
099 coordination mechanisms (Figure 1). As illustrated 151  
100 in Figure 1(a), MAPGD decomposes prompt op- 152  
101 timization into multiple orthogonal refinement di- 153  
102 mensions and assigns each to a specialized agent. 154  
103 Each agent independently analyzes model errors 155  
104 and produces textual pseudo-gradients that capture 156  
105 targeted improvement directions. 157

106 However, pseudo-gradients produced by dif- 158  
107 ferent agents may encode semantically conflict- 159  
108 ing refinement intents. To explicitly disentangle 160  
109 such conflicts, MAPGD proposes Hypersphere- 161  
110 Constrained Gradient Clustering (HCGC), which 162  
111 embeds agent-generated gradients into a shared se- 163  
112 mantic space and separates incompatible optimiza- 164  
113 tion directions via angularly constrained clustering. 165  
114 This design promotes semantic consistency within 166  
115 clusters while preventing destructive interference 167  
116 across conflicting refinement signals. 168

117 While HCGC addresses conflicts at the direc- 169  
118 tional level, agent contributions within the same se- 170  
119 mantic cluster may still vary in reliability. MAPGD 171  
120 therefore incorporates CAAW, which dynamically 172  
121 reweights agent channels based on validation per- 173  
122 formance. By emphasizing consistently effec- 174  
123 tive signals and attenuating noisy or unstable 175  
124 ones, CAAW enables robust aggregation of non- 176  
125 conflicting gradients, leading to stable and effective 177  
126 prompt refinement. 178

127 In summary, our contributions are three-fold. 179

128 (1) We propose MAPGD, a multi-agent pseudo- 180  
129 gradient framework for stable and interpretable 181  
130 prompt optimization. 182

131 (2) We introduce HCGC and CAAW to disentangle 183  
132 conflicting refinement signals and adaptively 184  
133 aggregate reliable agent contributions. 185

134 (3) We establish theoretical convergence guar- 186  
135 antees, showing that MAPGD preserves classical 187

136 gradient descent behavior with almost sure conver- 137  
138 gence to a local optimum at rate  $\mathcal{O}(1/\sqrt{T})$ . 139

## 2 Related Work 138

**Prompt Learning and Optimization.** Early 139  
140 prompt engineering relied on manual design, which 141  
142 is labor-intensive and difficult to scale. Automated 143  
144 approaches have since been explored, including re- 145  
146 inforcement learning (Deng et al., 2022), evolution- 147  
148 ary search (Fernando et al., 2023), and Bayesian 149  
150 optimization (Sahoo et al., 2025). Continuous 151  
152 prompt tuning methods, such as prefix tuning (Li 153  
154 and Liang, 2021) and soft prompts (Lester et al., 155  
156 2021), optimize prompts in embedding space but 157  
158 often sacrifice interpretability and require access 159  
160 to model internals. MAPGD operates in natural 161  
162 language space, preserving interpretability while 163  
164 enabling feedback-driven optimization. 165

**Gradient-Inspired Prompt Optimization.** Re- 166  
167 cent work has framed prompt optimization as a 168  
169 gradient-like process in discrete language space. 170  
171 ProTeGi (Pryzant et al., 2023) generates natural 172  
173 language pseudo-gradients from model errors and 174  
175 iteratively refines prompts via beam search and 176  
177 bandit-based selection. However, as a single-agent 178  
179 method, ProTeGi entangles heterogeneous refine- 180  
181 ment objectives and lacks explicit mechanisms for 182  
183 resolving conflicting optimization signals, which 184  
185 can lead to unstable or suboptimal convergence. 186  
187 This limitation motivates approaches that can dis- 188  
189 entangle and coordinate diverse refinement direc- 189  
190 tions. 190

**Multi-Agent Collaboration.** Multi-agent formu- 167  
168 lations enable parallel exploration of diverse per- 168  
169 spectives and have been widely studied in rein- 169  
170 forcement learning and distributed AI. In NLP, 170  
171 multi-agent debate (Liang et al., 2023; Du et al., 171  
172 2023) and collaborative generation (Hong et al., 172  
173 2024) demonstrate the benefits of role specializa- 173  
174 tion. However, increased diversity makes semantic 174  
175 conflicts inevitable, as different agents may pro- 175  
176 pose incompatible updates. MAPGD adopts a 176  
177 multi-agent design with specialized roles, while 177  
178 explicitly addressing the resulting coordination and 178  
179 conflict resolution challenges. 179

**Geometric and Adaptive Coordination.** Re- 180  
181 solving semantic conflicts requires structured in- 181  
182 ductive biases beyond heuristic fusion. In repre- 182  
183 sentation learning, contrastive objectives and an- 183  
184 gular margin constraints enforce intra-cluster com- 184

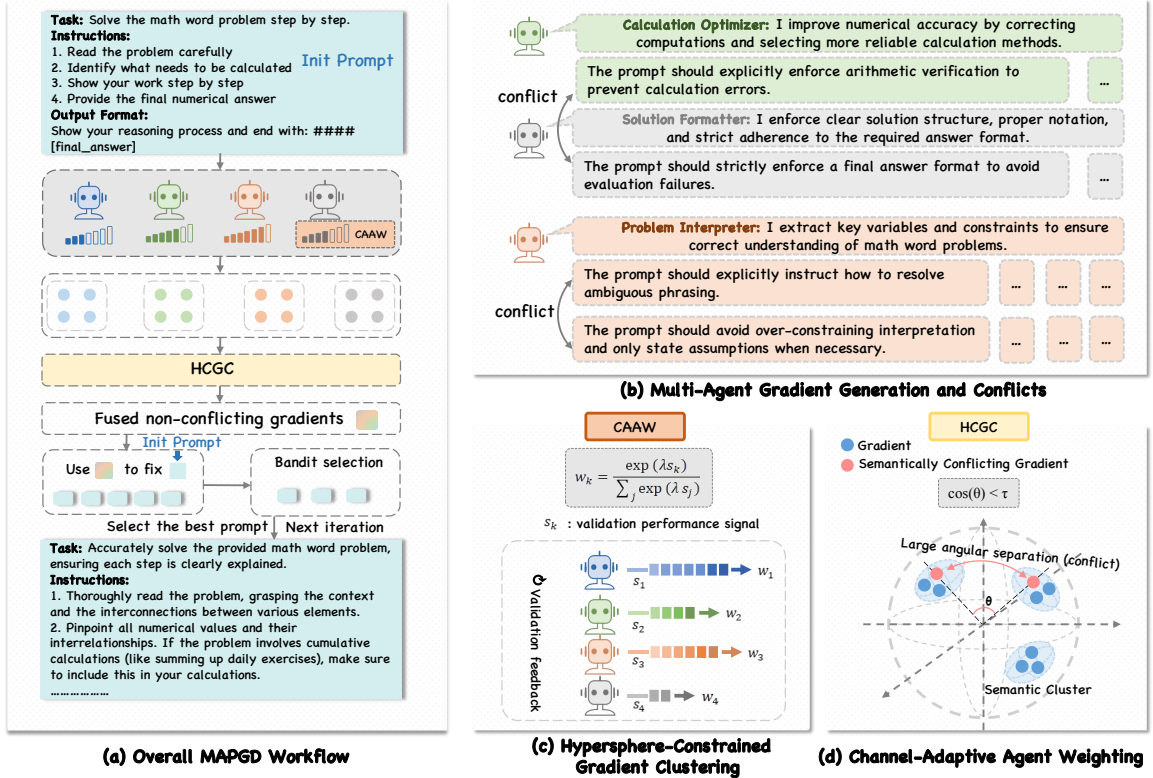


Figure 1: Overview of MAPGD. MAPGD performs prompt optimization through collaborative multi-agent pseudo-gradients. Agent-generated gradients are embedded and clustered via HCGC to resolve semantic conflicts, and adaptively weighted by CAAW according to validation performance. The resulting fused gradients iteratively refine the prompt through candidate generation and selection.

185 pactness and inter-cluster separation (Wang and  
 186 Isola, 2020; Liu et al., 2017; Wang et al., 2018;  
 187 Deng et al., 2019). Chen et al. (Chen et al., 2021)  
 188 further enhance large-margin contrastive learning  
 189 via distance polarization regularization. Inspired  
 190 by these advances, MAPGD introduces HCGC  
 191 for geometry-aware gradient coordination, com-  
 192 plemented by CAAW to adaptively calibrate agent  
 193 contributions.

194 **Concurrent Work.** Concurrent studies such as  
 195 MAPRO (Zhang et al., 2025) and MA-SAPO (Seo  
 196 et al., 2025) also explore multi-agent prompt op-  
 197 timization under different assumptions, including  
 198 explicit interaction structures or external memory  
 199 components. These settings are not directly aligned  
 200 with our budgeted, task-level optimization protocol.  
 201 A detailed comparison is provided in Appendix E.

202 **In summary,** MAPGD unifies prompt optimiza-  
 203 tion, gradient-inspired refinement, multi-agent col-  
 204 laboration, and geometry-aware coordination into  
 205 an interpretable framework that explicitly addresses  
 206 semantic conflicts.

### 207 3 Methodology

#### 208 3.1 Framework Overview

209 MAPGD formulates prompt optimization as a hy-  
 210 brid discrete–continuous gradient descent process  
 211 operating directly in natural language space. Un-  
 212 like continuous embedding-based approaches such  
 213 as prefix tuning or soft prompt optimization (Li  
 214 and Liang, 2021; Lester et al., 2021), MAPGD  
 215 preserves interpretability by explicitly manipulat-  
 216 ing textual prompts while still leveraging gradient-  
 217 inspired feedback signals for iterative refinement.  
 218 As illustrated in Figure 1, MAPGD frames prompt  
 219 optimization as a collaborative multi-agent process,  
 220 in which specialized agents focus on orthogonal  
 221 refinement dimensions and produce textual pseudo-  
 222 gradients that guide prompt updates.

223 Building on this formulation, MAPGD instanti-  
 224 ates the optimization procedure through an iterative  
 225 coordination loop. At each iteration, specialized  
 226 agents independently examine the current prompt  
 227 and available task feedback, and generate natural  
 228 language pseudo-gradients that reflect distinct re-  
 229 finement perspectives. These gradients are embed-

ded into a shared semantic space, where potential conflicts are explicitly identified and mitigated via geometry-aware clustering and adaptive fusion. The resulting fused gradients guide the generation of successor prompts, which are subsequently evaluated and filtered under a budgeted selection mechanism. Repeating this process yields an evolving optimization trajectory that systematically balances exploration across heterogeneous refinement directions with exploitation of effective updates.

Formally, the prompt optimization objective is defined as

$$F(p) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(M(x;p), y)], \quad (1)$$

where  $\ell(\cdot, \cdot)$  denotes a task-specific loss function and  $M(x;p)$  represents the model output conditioned on prompt  $p$ . The optimization goal is

$$p^* = \arg \min_p F(p). \quad (2)$$

Since gradients in this discrete textual space are not directly accessible, MAPGD constructs natural language pseudo-gradients,

$$\nabla F(p^{(t)}) \approx g^{(t)}, \quad (3)$$

which serve as semantic analogues of numerical gradients in classical optimization.

### 3.2 Specialized Prompt Agents

MAPGD employs multiple specialized agents, each responsible for a distinct and approximately orthogonal dimension of prompt refinement. Typical agent roles include instruction clarity, example selection, output formatting, and stylistic refinement, optionally complemented by a generic agent for broad exploratory updates. At iteration  $t$ , the resulting pseudo-gradient set is

$$G^{(t)} = \{g_1^{(t)}, g_2^{(t)}, \dots, g_K^{(t)}\}.$$

By decomposing prompt optimization into specialized perspectives, MAPGD enables parallel exploration of heterogeneous refinement directions, alleviating the limited signal diversity and premature convergence commonly observed in single-agent optimization.

### 3.3 Hypersphere-Constrained Gradient Clustering

Inspired by hyperspherical metric learning (Wang et al., 2018; Liu et al., 2017), HCGC resolves

semantic conflicts among heterogeneous pseudo-gradients by enforcing geometry-aware separation on a normalized hypersphere. Instead of assuming all agent proposals are mutually compatible, HCGC explicitly detects semantic disagreement and clusters gradients into coherent refinement directions before fusion.

#### Gradient Embedding and Conflict Signal.

Given agent-generated pseudo-gradients  $G^{(t)}$ , each  $g_k^{(t)}$  is encoded into a  $d$ -dimensional semantic vector  $v_k^{(t)} = \phi(g_k^{(t)})$  and normalized onto the unit hypersphere:

$$\hat{v}_k^{(t)} = \frac{v_k^{(t)}}{\|v_k^{(t)}\|}, \quad \hat{v}_k^{(t)} \in \mathbb{S}^{d-1}. \quad (4)$$

Semantic similarity is measured by cosine similarity,

$$\text{sim}(\hat{v}_i, \hat{v}_j) = \hat{v}_i^\top \hat{v}_j = \cos(\Delta(\hat{v}_i, \hat{v}_j)), \quad (5)$$

where  $\Delta(\cdot, \cdot)$  denotes the angular distance on the hypersphere. To make ‘‘conflict’’ explicit and measurable, we define two gradients as being in semantic conflict if their similarity falls below a threshold:

$$\text{conflict}(\hat{v}_i, \hat{v}_j) \Leftrightarrow \hat{v}_i^\top \hat{v}_j < \theta_{\text{conflict}}. \quad (6)$$

Intuitively, low cosine similarity indicates substantially different (and potentially incompatible) refinement intents, motivating explicit separation in clustering.

#### Clustering and Adaptive Cluster Cardinality.

We perform cosine K-means on the hypersphere to group semantically coherent directions. Rather than fixing the number of clusters, we adopt a bounded adaptive strategy:

$$K^{(t)} = \min(|G^{(t)}|, K_{\text{max}}), \quad (7)$$

where  $|G^{(t)}|$  is the number of gradients at iteration  $t$  and  $K_{\text{max}}$  is an upper bound. This prevents over-fragmentation when gradients are few, while allowing multiple distinct directions when agent proposals are diverse.

**Angular Margin Constraint.** Initial clustering may still yield weakly separated clusters when conflicting gradients coexist. HCGC therefore refines assignments with an angular margin constraint. Let  $u_i$  be the centroid of  $\hat{v}_k$ ’s assigned cluster and  $u_j$  any other centroid. Define  $\alpha = \Delta(\hat{v}_k, u_i)$  and

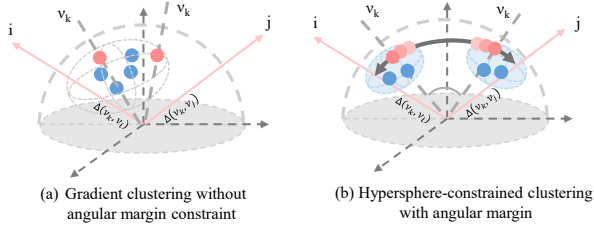


Figure 2: Effect of hypersphere-constrained angular margin in HCGC. (a) Clustering without angular margin constraints may entangle semantically conflicting pseudo-gradients, leading to ambiguous cluster boundaries. (b) Enforcing an angular margin on the unit hypersphere improves intra-cluster compactness and inter-cluster separation, effectively isolating conflicting refinement directions.

$\beta = \Delta(\hat{v}_k, u_j)$ . While standard assignment requires  $\alpha < \beta$ , HCGC strengthens it by a margin scale  $n \geq 1$ :

$$n \cdot \alpha < \beta, \quad (8)$$

equivalently,

$$\cos(n \cdot \alpha) > \cos(\beta). \quad (9)$$

Gradients that violate the margin are reassigned to clusters where the constraint holds, yielding more compact and better-separated semantic groups.

**Fusion of Clustered Gradients.** Finally, gradients within each cluster are fused into a single representative refinement direction via LLM-guided synthesis. To avoid redundant updates, we apply a diversity filter to discard fused gradients overly similar to those already selected, controlled by  $\theta_{\text{diversity}}$ .

### 3.4 Channel-Adaptive Agent Weighting

Although HCGC resolves structural conflicts between semantic gradient directions, gradients within the same cluster may still vary in reliability due to heterogeneous agent expertise and historical effectiveness. To account for this variability, we introduce CAAW, a reliability-aware coordination mechanism that adaptively calibrates agent contributions during gradient fusion. CAAW is inspired by channel-wise recalibration mechanisms in neural representations (Hu et al., 2018), but operates at the level of agent-generated textual gradients.

**Adaptive Weight Estimation.** Let  $s_k^{(t)}$  denote the validation performance gain attributed to gradient  $g_k^{(t)}$  at iteration  $t$ , estimated from historical agent performance statistics. CAAW converts these

scalar signals into normalized reliability weights via a softmax transformation:

$$w_k^{(t)} = \frac{\exp(\lambda s_k^{(t)})}{\sum_{j \in C} \exp(\lambda s_j^{(t)})}, \quad (10)$$

where  $C$  denotes the set of gradients within the same semantic cluster and  $\lambda$  controls the sharpness of the weighting distribution. This formulation assigns higher weights to agents that have consistently produced effective refinements, while down-weighting less reliable channels.

**Weight-Conditioned LLM Fusion.** Rather than explicitly performing an arithmetic weighted sum over textual gradients, CAAW injects the estimated weights into the LLM-driven fusion operator as conditioning signals. Concretely, the fusion operator

$$g_{\text{fused}}^{(t)} = \Psi(\{g_k^{(t)}\}_{k \in C} \mid \{w_k^{(t)}\}_{k \in C}) \quad (11)$$

receives both the set of clustered gradients and their corresponding reliability weights, which are encoded in the fusion prompt to bias the synthesis process toward more trustworthy agent suggestions. In this way, agent weights influence the fused gradient implicitly through prompt-level emphasis, rather than through explicit numerical combination.

**Stability Implications.** By conditioning gradient fusion on adaptive reliability weights, CAAW suppresses noisy or redundant updates while amplifying consistently effective refinement channels. This weight-aware fusion mechanism reduces variance in the optimization trajectory and complements HCGC by addressing intra-cluster reliability differences, thereby promoting stable convergence toward semantically robust prompt candidates.

### 3.5 Candidate Generation and Selection

Fused gradients generate successor prompts  $\{p'_1, \dots, p'_n\}$ . MAPGD applies lightweight filtering and optional bandit-based evaluation (Audibert and Bubeck, 2010) to keep the computation budgeted. Unlike prior work that depends heavily on bandit selection, MAPGD attains robustness and interpretability mainly via conflict-aware clustering and adaptive weighting across agents. Overall algorithm can be found in Algorithm 1.

### 3.6 Theoretical Convergence

We provide a theoretical analysis of MAPGD in appendix G. Here, we summarize the key assumptions

---

**Algorithm 1** MAPGD with HCGC and CAAW

---

**Require:** Initial prompt  $p_0$ , train data  $D_{\text{train}}$ , dev data  $D_{\text{dev}}$ , agents  $\{A_i\}_{i=1}^N$ , iterations  $R$ , beam width  $k$

- 1: Initialize best prompt  $p_\star^{(0)} \leftarrow p_0$ , beam  $B_0 \leftarrow \{p_0\}$ ; initialize each agent prompt  $A_i.p \leftarrow p_0$
  - 2: **for**  $t = 1$  to  $R$  **do**
  - 3:    $M_t \leftarrow \text{SAMPLEMINIBATCH}(D_{\text{train}}, b)$
  - 4:    $\mathcal{G}_t \leftarrow \text{GENEAGENTGRADS}(\{A_i\}, M_t) \triangleright$   
    Alg. 2
  - 5:    $\tilde{\mathcal{G}}_t \leftarrow \text{HCGC}(\mathcal{G}_t) \quad \triangleright$  Alg. 3
  - 6:    $F_t \leftarrow \text{CAAW}(\tilde{\mathcal{G}}_t, D_{\text{dev}}) \quad \triangleright$  Alg. 4
  - 7:    $\mathcal{C}_t \leftarrow \text{EXPANDPROMPTS}(p_\star^{(t-1)}, F_t)$
  - 8:    $B_t, p_\star^{(t)} \leftarrow \text{BANDITSELECT}(\mathcal{C}_t, D_{\text{dev}}, k)$
  - 9:    $\text{SYNCAGENTS}(\{A_i\}, p_\star^{(t)}) \quad \triangleright$  Alg. 5
  - 10:   **if**  $\text{CONVERGED}(p_\star^{(t)})$  **then break**
  - 11:   **end if**
  - 12: **end for**
  - 13: **return** Best prompt over  $\{p_\star^{(1)}, \dots, p_\star^{(t)}\}$
- 

and main implications to clarify how the proposed framework relates to classical stochastic optimization.

Specifically, our analysis relies on three standard conditions: (i) bounded second moments of agent-generated pseudo-gradients, (ii) smoothness (or Lipschitz continuity) of the empirical loss function in the semantic embedding space, (iii) unbiased sampling induced by the bandit-based candidate selection mechanism. Under these assumptions, MAPGD admits an interpretation as a stochastic approximation process and achieves a sublinear convergence rate of  $\mathcal{O}(1/\sqrt{T})$  in both convex and non-convex settings.

Intuitively, the core components of MAPGD play complementary roles in ensuring stable optimization. HCGC aggregates semantically aligned refinement directions and mitigates conflicting updates, reducing variance in the fused pseudo-gradient. CAAW further emphasizes agents that have demonstrated consistent utility across iterations. Finally, bandit-based candidate selection maintains sufficient exploration while preventing systematic bias. Together, these mechanisms ensure that the overall optimization dynamics remain stable and progressively improve prompt quality over iterations.

## 4 Experiments

Experiments evaluate MAPGD in terms of effectiveness, efficiency, and robustness. For fair comparison, we reproduce ProTeGi under identical set-

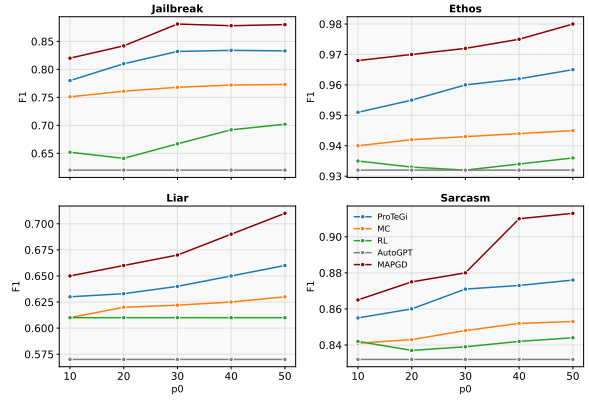


Figure 3: Test performance (F1 score) vs. API query budget per prompt candidate across four benchmark tasks (evaluated up to 50 queries for fair comparison with ProTeGi).

tings and further conduct ablations and token consumption analysis. A real-world text-generation case study is also included (Appendix I) to demonstrate applicability beyond benchmarks.

### 4.1 Experimental Details

**Datasets.** Following ProTeGi’s protocol, we evaluate on four classification benchmarks (Jailbreak, Ethos, LIAR, Sarcasm) and three arithmetic reasoning datasets (GSM8K, AQUARAT, SVAMP). For each dataset, we sample 50/150 instances for dev/test and report test F1 (classification) or accuracy (reasoning). Dataset details are in Appendix A.

**Setup.** We use the January 2023 gpt-3.5-turbo via Azure OpenAI unless stated otherwise (temperature 0.0 for classification). Results are averaged over three runs; F1 is computed by max-pooling over the final prompt beam. All methods share identical initial prompts, data, and random seeds, with a fixed budget of  $T=10$  and default hyperparameters. Unless stated otherwise, MAPGD follows Section 3.3 with  $\phi=\text{all-MiniLM-L6-v2}$ ,  $\theta_{\text{conflict}}=0.3$ , clustering threshold 0.7,  $\theta_{\text{diversity}}=0.7$ , and UCB selection with 80 evaluations;  $\lambda=1$  (Appendix D). All experiments in Section 4.2 use these modules unless ablated.

**Baselines.** We primarily compare against ProTeGi (Pryzant et al., 2023) under identical query budgets and initialization, and additionally include MC search, RL-style edit methods, AutoGPT, and three reasoning baselines (InsZero, Instinct, PromptWizard). Full descriptions are in Appendix B.

### 4.2 Experimental Results

**Classification benchmarks.** Figure 3 reports test F1 under different per-candidate evaluation bud-

Method	GSM8k	AQUARAT	SVAMP
InsZero	74.2	54.3	79.5
Instinct	74.5	54.7	81.0
PromptWizard	90.0	58.2	82.3
MAPGD	<b>93.5</b>	<b>60.3</b>	<b>84.1</b>

Table 1: Arithmetic reasoning accuracy (%) on GSM8k, AQUARAT, and SVAMP. Best results in bold.

gets (up to 50 queries, following ProTeGi for fairness). MAPGD consistently achieves the best performance across all four datasets, with most gains emerging within the first 30–40 evaluations. The improvements are most pronounced on conflict-prone tasks, suggesting that *explicit conflict disentanglement and adaptive coordination* are critical when refinement signals disagree: on **Jailbreak**, MAPGD reaches **0.880** at 50 queries (vs. 0.833 for ProTeGi), and on **LIAR** it attains **0.710** (vs. 0.660). We attribute these gains to HCGC, which clusters compatible pseudo-gradients before fusion to reduce semantic interference, and CAAW, which down-weights unreliable agents to prevent noisy edits from dominating early iterations, thereby improving budget efficiency. On the more saturated **Ethos** dataset, MAPGD still improves to **0.980** (vs. 0.965), indicating that coordinated multi-agent feedback can extract residual refinements even near a performance ceiling. We follow the 50-query budget for strict comparability; extended-budget results are deferred to Appendix C.

**Arithmetic reasoning benchmarks.** Table 1 evaluates MAPGD on GSM8k, AQUARAT, and SVAMP, compared against InsZero, Instinct, and PromptWizard. MAPGD achieves consistent gains across all three datasets, most notably on **GSM8k (93.5%)**, improving over the strongest baseline by +3.5 points. We hypothesize that the benefit on reasoning stems from the same coordination principle: specialized agents propose complementary edits (e.g., step-by-step constraints, format strictness, and exemplar guidance), while HCGC/CAAW suppress incompatible or low-yield directions, leading to more stable chain-of-thought prompting. MAPGD also improves the best baseline on **AQUARAT (60.3% vs. 58.2%)** and **SVAMP (84.1% vs. 82.3%)**, confirming transfer beyond classification.

### 4.3 Ablation Study

To isolate the sources of MAPGD’s gains, we conduct a component-wise ablation by selectively disabling HCGC and CAAW. Table 2 reports four

Table 2: Ablation study on core MAPGD components. ✓ indicates the component is enabled. Results are F1 score (mean over 3 runs).

Setting	HCGC	CAAW	Jailbreak	Ethos	LIAR	Sarcasm
I	×	×	0.82	0.91	0.62	0.83
II	×	✓	0.84	0.94	0.62	0.87
III	✓	×	0.85	0.95	0.63	0.88
IV	✓	✓	<b>0.88</b>	<b>0.98</b>	<b>0.65</b>	<b>0.91</b>

settings (I–IV), where Setting I removes both modules and serves as a minimal baseline. Enabling either module leads to consistent improvements. Introducing HCGC (III vs. I) yields systematic gains, confirming that geometry-aware clustering prior to fusion is essential for reducing semantic interference among pseudo-gradients. CAAW alone (II vs. I) also improves performance, indicating that uniform averaging is suboptimal under varying agent reliability; adaptive reweighting based on validation feedback enhances robustness. The combined setting (IV) consistently outperforms either component alone, suggesting strong complementarity: HCGC resolves directional incompatibility via conflict disentanglement, whereas CAAW addresses reliability mismatch through adaptive credit assignment. Gains are smaller on saturated datasets (e.g., Ethos), consistent with ceiling effects, while conflict-prone tasks benefit more from structured coordination.

### 4.4 Robustness to Dataset Sampling Size

Following the standard ProTeGi evaluation protocol, our main results use 50 development and 150 test examples per dataset to ensure direct comparability. To verify that MAPGD’s gains are not an artifact of small evaluation sets, we additionally evaluate larger subsamples on three representative classification benchmarks.

Table 3 reports results under three sampling regimes (50/150, 100/300, and 200/500). MAPGD consistently outperforms ProTeGi and MC across all datasets and sampling sizes, indicating that its improvements persist as the evaluation set grows.

### 4.5 Token Consumption Analysis

Token efficiency is evaluated under the setup of Section 4.1 by counting both input and output tokens per API call. We report total tokens, number of calls, and performance-per-token. As shown in Table 4, MAPGD attains higher F1 with fewer tokens and calls than ProTeGi, reducing per-iteration token cost by ~8% and improving performance-per-token by >10%, suggesting reduced redundant

Dataset	Method	50/150	100/300	200/500
Jailbreak	ProTeGi	0.82	0.83	0.81
	MC	0.76	0.77	0.75
	<b>MAPGD</b>	<b>0.86</b>	<b>0.86</b>	<b>0.84</b>
Ethos	ProTeGi	0.93	0.94	0.94
	MC	0.89	0.90	0.90
	<b>MAPGD</b>	<b>0.98</b>	<b>0.97</b>	<b>0.97</b>
LIAR	ProTeGi	0.65	0.62	0.63
	MC	0.62	0.63	0.63
	<b>MAPGD</b>	<b>0.71</b>	<b>0.71</b>	<b>0.69</b>

Table 3: Robustness to dataset sampling size (F1 score) under three dev/test splits.

Metric	ProTeGi	MAPGD (full)
Total Tokens	256k	<b>236k</b>
Calls	962	<b>643</b>
Avg F1	0.83	<b>0.87</b>
F1 / Token ( $\times 10^{-6}$ )	3.24	<b>3.69</b>

Table 4: Token consumption comparison between ProTeGi and MAPGD (averaged over 3 runs). MAPGD achieves lower token usage and higher efficiency.

Component	ProTeGi (s)	MAPGD (s)
Gradient generation	10.2	37.9
Embedding + HCGC	–	24.4
Prompt expansion	33.7	28.7
Bandit selection	30.3	26.5
Other	85.2	84.0
<b>Total</b>	<b>159.4</b>	<b>201.5</b>

Table 5: Runtime breakdown (in seconds) comparing ProTeGi and MAPGD.

queries via HCGC and CAAW.

#### 4.6 Computational Overhead Analysis

We compare runtime overhead against ProTeGi under identical settings. MAPGD shares evaluation, prompt expansion, and bandit selection with ProTeGi, but adds multi-agent gradient generation and embedding-based clustering. As shown in Table 5, the extra overhead mainly comes from multi-agent gradient generation, which scales linearly with the number of agents. Embedding and HCGC incur an almost constant cost, and with four agents the overall overhead remains a small fraction of total runtime.

#### 4.7 Agent Number Sensitivity

To motivate the use of four specialized agents, we vary the number of agents  $N \in \{2, 4, 6\}$  under the setup of Section 4.1. For  $N = 2$ , only the Instruction Specialist and Example Curator are used; for  $N = 6$ , two additional style/format agents are added. As shown in Figure 4, average F1 improves

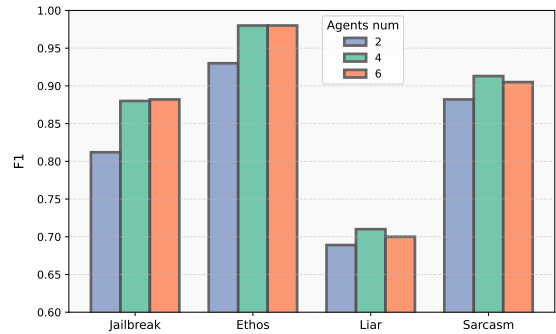


Figure 4: Agent number sensitivity analysis. Increasing agents from 2 to 4 yields clear improvements, while further increasing to 6 brings diminishing returns.

substantially from  $N = 2$  to  $N = 4$ , while  $N = 6$  yields only marginal gains with higher token cost, indicating diminishing returns. Thus,  $N = 4$  offers the best effectiveness–efficiency trade-off.

#### 4.8 Experimental Insights

MAPGD consistently outperforms single-agent ProTeGi and other baselines, with the largest gains on LIAR and Jailbreak where refinement signals often conflict. Ablations show HCGC and CAAW are complementary: removing either hurts performance, and removing both approaches the single-agent regime. MAPGD is also more token-efficient than ProTeGi (fewer API calls,  $>10\%$  higher performance-per-token) and transfers well to arithmetic reasoning, surpassing strong baselines on GSM8k, AQUARAT, and SVAMP. Appendix H provides a qualitative case study illustrating how HCGC and CAAW resolve conflicting pseudo-gradients.

### 5 Conclusion

This work introduces MAPGD, a coordinated multi-agent framework for prompt optimization that integrates HCGC and CAAW to resolve conflicting refinement directions and reduce the impact of unreliable agents. Across classification and arithmetic reasoning benchmarks, MAPGD consistently improves over strong baselines while lowering token consumption, and it admits a sublinear convergence guarantee of  $O(1/\sqrt{T})$  under standard stochastic approximation assumptions (Appendix G). Overall, MAPGD provides an interpretable, robust, and resource-efficient approach to prompt optimization via structured multi-agent collaboration.

## 595 Limitations

596 MAPGD relies on semantic embeddings to detect  
597 conflicts and cluster pseudo-gradients; under do-  
598 main shift or highly specialized jargon, similarity  
599 estimates may become unreliable and degrade co-  
600 ordination quality. The method also depends on  
601 reasonably informative agent feedback: if multi-  
602 ple agents are systematically misaligned, adaptive  
603 weighting may not fully suppress harmful edits. Fi-  
604 nally, compared to single-agent methods, MAPGD  
605 introduces extra overhead from multi-agent gen-  
606 eration and coordination (embedding/clustering),  
607 which may limit deployment under strict latency or  
608 budget constraints.

## 609 References

610 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
611 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
612 Diogo Almeida, Janko Altenschmidt, Sam Altman,  
613 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-  
614 cal report. *arXiv preprint arXiv:2303.08774*.

615 Eshaan Agarwal, Joykirat Singh, Vivek Dani, Raghav  
616 Magazine, Tanuja Ganu, and Akshay Nambi. 2024.  
617 [Promptwizard: Task-aware prompt optimization](#)  
618 [framework](#). *Preprint*, arXiv:2405.18369.

619 Jean-Yves Audibert and Sébastien Bubeck. 2010. Best  
620 arm identification in multi-armed bandits. In *COLT-*  
621 *23th Conference on learning theory-2010*, pages 13–  
622 p.

623 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan,  
624 Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter  
625 Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, and  
626 1 others. 2023. Sparks of artificial general intelli-  
627 gence: Early experiments with gpt-4. *arXiv preprint*  
628 *arXiv:2303.12712*.

629 Lichang Chen, Jiu Hai Chen, Tom Goldstein, Heng  
630 Huang, and Tianyi Zhou. 2023. [Instructzero: Ef-](#)  
631 [ficient instruction optimization for black-box large](#)  
632 [language models](#). *Preprint*, arXiv:2306.03082.

633 Shuo Chen, Gang Niu, Chen Gong, Jun Li, Jian Yang,  
634 and Masashi Sugiyama. 2021. Large-margin con-  
635 trastive learning with distance polarization regular-  
636 izer. In *International conference on machine learn-*  
637 *ing*, pages 1673–1683. PMLR.

638 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,  
639 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias  
640 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro  
641 Nakano, and 1 others. 2021. Training verifiers  
642 to solve math word problems. *arXiv preprint*  
643 *arXiv:2110.14168*.

644 Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos  
645 Zafeiriou. 2019. Arcface: Additive angular margin  
646 loss for deep face recognition. In *Proceedings of*

*the IEEE/CVF conference on computer vision and*  
*pattern recognition*, pages 4690–4699. 647  
648

Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yihan  
Wang, Han Guo, Tianmin Shu, Meng Song, Eric P  
Xing, and Zhiting Hu. 2022. Rlprompt: Optimizing  
discrete text prompts with reinforcement learning.  
*arXiv preprint arXiv:2205.12548*. 649  
650  
651  
652  
653

Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenen-  
baum, and Igor Mordatch. 2023. Improving factual-  
ity and reasoning in language models through multi-  
agent debate. In *Forty-first International Conference*  
*on Machine Learning*. 654  
655  
656  
657  
658

Ibrahim Abu Farha and Walid Magdy. 2020. From  
arabic sentiment analysis to sarcasm detection: The  
arsarcasm dataset. In *The 4th Workshop on Open-*  
*Source Arabic Corpora and Processing Tools*, pages  
32–39. European Language Resources Association  
(ELRA). 659  
660  
661  
662  
663  
664

Chrisantha Fernando, Dylan Banarse, Charles Blun-  
dell, Tim Rocktäschel, and Simon Osindero.  
2023. Promptbreeder: Self-referential self-  
improvement via prompt evolution. *arXiv preprint*  
*arXiv:2309.16797*. 665  
666  
667  
668  
669

Sirui Hong, Mingchen Zhuge, Jiaqi Chen, Xiawu Zheng,  
Yuheng Cheng, Ceyao Zhang, Jinlin Wang, Zili  
Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang  
Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu,  
and Jürgen Schmidhuber. 2024. [Metagtpt: Meta pro-](#)  
[gramming for a multi-agent collaborative framework](#).  
*Preprint*, arXiv:2308.00352. 670  
671  
672  
673  
674  
675  
676

Jie Hu, Li Shen, and Gang Sun. 2018. Squeeze-and-  
excitation networks. In *Proceedings of the IEEE con-*  
*ference on computer vision and pattern recognition*,  
pages 7132–7141. 677  
678  
679  
680

Ellen Jiang, Kristen Olson, Edwin Toh, Alejandra  
Molina, Aaron Donsbach, Michael Terry, and Carrie J  
Cai. 2022. Promptmaker: Prompt-based prototyping  
with large language models. In *CHI Conference on*  
*Human Factors in Computing Systems Extended Ab-*  
*stracts*, pages 1–8. 681  
682  
683  
684  
685  
686

Brian Lester, Rami Al-Rfou, and Noah Constant. 2021.  
The power of scale for parameter-efficient prompt  
tuning. *arXiv preprint arXiv:2104.08691*. 687  
688  
689

Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning:](#)  
[Optimizing continuous prompts for generation](#). *Pro-*  
*ceedings of the 59th Annual Meeting of the Asso-*  
*ciation for Computational Linguistics and the 11th*  
*International Joint Conference on Natural Language*  
*Processing (Volume 1: Long Papers)*, pages 4582–  
4597. 690  
691  
692  
693  
694  
695  
696

Weizhe Liang, Yizhong Wu, Mina Lee, Ed Chi, Denny  
Zhou, Yiming Song, and Xiang Lisa Li. 2023. En-  
couraging divergent thinking in large language mod-  
els through multi-agent debate. *arXiv preprint*  
*arXiv:2305.19118*. 697  
698  
699  
700  
701

702	Xiaoqiang Lin, Zhaoxuan Wu, Zhongxiang Dai,	Taylor Shin, Yasaman Razeghi, Robert L Logan IV,	757
703	Wenyang Hu, Yao Shu, See-Kiong Ng, Patrick Jail-	Eric Wallace, and Sameer Singh. 2020. Autoprompt:	758
704	let, and Bryan Kian Hsiang Low. 2024. <a href="#">Use your</a>	Eliciting knowledge from language models with	759
705	<a href="#">instinct: Instruction optimization for llms using neu-</a>	automatically generated prompts. <i>arXiv preprint</i>	760
706	<a href="#">ral bandits coupled with transformers.</a> <i>Preprint,</i>	<i>arXiv:2010.15980.</i>	761
707	arXiv:2310.02905.		
708	Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blun-	Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Di-	762
709	som. 2017. Program induction by rationale genera-	hong Gong, Jingchao Zhou, Zhifeng Li, and Wei Liu.	763
710	tion: Learning to solve and explain algebraic word	2018. Cosface: Large margin cosine loss for deep	764
711	problems. <i>arXiv preprint arXiv:1705.04146.</i>	face recognition. In <i>Proceedings of the IEEE con-</i>	765
712	Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhik-	<i>ference on computer vision and pattern recognition,</i>	766
713	sha Raj, and Le Song. 2017. Spheroface: Deep hy-	pages 5265–5274.	767
714	persphere embedding for face recognition. In <i>Pro-</i>	Tongzhou Wang and Phillip Isola. 2020. Understanding	768
715	<i>ceedings of the IEEE conference on computer vision</i>	contrastive representation learning through alignment	769
716	<i>and pattern recognition,</i> pages 212–220.	and uniformity on the hypersphere. In <i>International</i>	770
717	Ioannis Mollas, Zoe Chrysopoulou, Stamatis Karlos,	<i>conference on machine learning,</i> pages 9929–9939.	771
718	and Grigorios Tsoumakas. 2022. Ethos: a multi-label	PMLR.	772
719	hate speech detection dataset. <i>Complex &amp; Intelligent</i>	W Wang. 2021. A new benchmark dataset for fake news	773
720	<i>Systems,</i> 8(6):4663–4678.	detection. In <i>Proceedings of the 55th Annual Meet-</i>	774
721	Arkil Patel, Satwik Bhattamishra, and Navin Goyal.	<i>ing of the Association for Computational Linguistics,</i>	775
722	2021. Are nlp models really able to solve	volume 2.	776
723	simple math word problems? <i>arXiv preprint</i>	J Diego Zamfirescu-Pereira, Richmond Y Wong, Bjoern	777
724	<i>arXiv:2103.07191.</i>	Hartmann, and Qian Yang. 2023. Why johnny can't	778
725	Archiki Prasad, Peter Hase, Xiang Zhou, and Mohit	prompt: how non-ai experts try (and fail) to design	779
726	Bansal. 2022. Grips: Gradient-free, edit-based in-	llm prompts. In <i>Proceedings of the 2023 CHI confer-</i>	780
727	struction search for prompting large language models.	<i>ence on human factors in computing systems,</i> pages	781
728	<i>arXiv preprint arXiv:2203.07281.</i>	1–21.	782
729	Reid Pryzant, Dan Iter, Jerry Li, Yin Tat Lee, Chen-	Tianjun Zhang, Xuezhi Wang, Denny Zhou, Dale	783
730	guang Zhu, and Michael Zeng. 2023. Automatic	Schuurmans, and Joseph E Gonzalez. 2022. Tem-	784
731	prompt optimization with "gradient descent" and	pera: Test-time prompting via reinforcement learning.	785
732	beam search. <i>arXiv preprint arXiv:2305.03495.</i>	<i>arXiv preprint arXiv:2211.11890.</i>	786
733	Guanghui Qin and Jason Eisner. 2021. Learning how	Zheyuan Zhang, Lin Ge, Hongjiang Li, Weicheng Zhu,	787
734	to ask: Querying lms with mixtures of soft prompts.	Chuxu Zhang, and Yanfang Ye. 2025. <a href="#">Mapro: Recast-</a>	788
735	<i>arXiv preprint arXiv:2104.06599.</i>	<a href="#">ing multi-agent prompt optimization as maximum a</a>	789
736	Laria Reynolds and Kyle McDonell. 2021. Prompt	<a href="#">posteriori inference.</a> <i>Preprint,</i> arXiv:2510.07475.	790
737	programming for large language models: Beyond	Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han,	791
738	the few-shot paradigm. In <i>Extended abstracts of the</i>	Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy	792
739	<i>2021 CHI conference on human factors in computing</i>	Ba. 2022. Large language models are human-level	793
740	<i>systems,</i> pages 1–7.	prompt engineers. In <i>The eleventh international con-</i>	794
741	Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha,	<i>ference on learning representations.</i>	795
742	Vinija Jain, Samrat Mondal, and Aman Chadha.	<b>A Datasets Details</b>	796
743	2025. <a href="#">A systematic survey of prompt engineering in</a>	To ensure comparability with ProTeGi, we evalu-	797
744	<a href="#">large language models: Techniques and applications.</a>	ate MAPGD on the same four benchmark NLP	798
745	<i>Preprint,</i> arXiv:2402.07927.	classification tasks that span diverse domains and	799
746	Wonduk Seo, Juhyeon Lee, Junseo Koh, Hyunjin	languages:	800
747	An, Jian Park, Seunghyun Lee, Haihua Chen, and	• <b>Jailbreak</b> (Shen et al., 2024): a novel task	801
748	Yi Bu. 2025. <a href="#">Prompt optimization via retrieved rea-</a>	to determine whether a user input to an LLM	802
749	<a href="#">soning assets and multi-agent analysis.</a> <i>Preprint,</i>	continuation API constitutes a jailbreak at-	803
750	arXiv:2510.16635.	tempt. Jailbreak attacks are defined as user	804
751	Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen,	strategies that aim to make the model violate	805
752	and Yang Zhang. 2024. "do anything now": Charac-	its own safety policies, such as generating	806
753	terizing and evaluating in-the-wild jailbreak prompts	harmful content or exposing hidden instruc-	807
754	on large language models. In <i>Proceedings of the</i>	tions. The dataset contains 1306 multilingual	808
755	<i>2024 on ACM SIGSAC Conference on Computer and</i>	examples with human-annotated jailbreak la-	809
756	<i>Communications Security,</i> pages 1671–1685.	breaks.	810

811	• <b>Ethos</b> (Mollas et al., 2022): an English hate speech detection dataset with 997 online comments annotated as hateful or non-hateful.	855
812		856
813		857
814	• <b>LIAR</b> (Wang, 2021): a large-scale dataset for fake news detection, consisting of 4,000 short statements labeled with ground-truth veracity, along with context information.	858
815		859
816		860
817		861
818	• <b>Sarcasm</b> (Farha and Magdy, 2020): an Arabic sarcasm detection dataset with 10,000 online comments, labeled for the presence or absence of sarcasm.	862
819		863
820		864
821		865
822	Beyond classification, we further evaluate on three arithmetic reasoning datasets that require multi-step logical inference and symbolic manipulation:	866
823		867
824		868
825		869
826	• <b>GSM8k</b> (Cobbe et al., 2021): a widely used benchmark of 8,500 grade-school math problems that require step-by-step reasoning.	870
827		871
828		872
829	• <b>AQUARAT</b> (Ling et al., 2017): a dataset of algebraic word problems designed to test symbolic reasoning and program induction.	873
830		874
831		875
832	• <b>SVAMP</b> (Patel et al., 2021): a benchmark focusing on simple arithmetic word problems with diverse linguistic variations, serving as a testbed for robustness to paraphrasing.	876
833		877
834		878
835		879
836	In addition, we include a real-world text generation case study in the Appendix to highlight MAPGD’s applicability beyond controlled benchmarks.	880
837		881
838		882
839		883
840	<b>B Baselines Details</b>	884
841	To evaluate the effectiveness of MAPGD, we compare it against a set of non-parametric prompt optimization methods, following the setup of ProTeGi (Pryzant et al., 2023), and additionally include ProTeGi itself as a strong baseline. Specifically, we consider:	885
842		886
843		887
844	• <b>ProTeGi</b> : the original prompt gradient descent framework, where a single agent iteratively generates pseudo-gradients and candidate prompts, with bandit-based selection. This serves as our primary baseline for comparison.	888
845		889
846		890
847	• <b>Monte-Carlo (MC)</b> (Zhou et al., 2022): an iterative but directionless Monte Carlo search	891
848		892
849		893
850		894
851		895
852		896
853		897
854		898
		899
	over the prompt space. For fairness, we match the number of samples per candidate to the successors generated by MAPGD.	855
		856
		857
	• <b>Reinforcement Learning (RL)</b> Recent approaches such as GrIPS (Prasad et al., 2022) and TEMPERA (Zhang et al., 2022) formulate prompt optimization as a reinforcement learning problem. In these methods, the prompt text is first segmented into phrases, and the search space is explored via phrase-level edit operations, including addition, paraphrasing, swapping, and deletion.	858
		859
		860
		861
		862
		863
		864
		865
		866
	• <b>AutoGPT</b> : an open-source autonomous agent system that improves prompts through self-directed feedback loops. We configure AutoGPT with the same number of examples and errors as MAPGD, running for the same number of optimization steps.	867
		868
		869
		870
		871
		872
	For arithmetic reasoning, we additionally compare against three competitive baselines, all reported with GPT-3.5-Turbo in a zero-shot setting:	873
		874
		875
	• <b>InsZero</b> (Chen et al., 2023): a zero-shot prompting method with instruction tuning.	876
		877
	• <b>Instinct</b> (Lin et al., 2024): an instruction-following baseline that improves prompting robustness.	878
		879
		880
	• <b>PromptWizard</b> (Agarwal et al., 2024): a state-of-the-art prompt synthesis framework that iteratively enriches prompts with intermediate reasoning steps and examples.	881
		882
		883
		884
	This collection of baselines allows us to evaluate MAPGD not only against single-agent gradient-based optimization, but also against strong arithmetic reasoning methods optimized for chain-of-thought style prompting.	885
		886
		887
		888
		889
	<b>C Extended Evaluation Budgets Beyond 50 Queries</b>	890
		891
	To verify whether our conclusions depend on the 50-query budget, we further extend the evaluation budget to <b>60</b> and <b>70</b> queries on three representative datasets (Jailbreak, LIAR, and Sarcasm). We keep the same evaluation protocol and report results in Table 6. Overall, the <b>method ranking remains unchanged</b> , and the additional gains beyond 50 are <b>marginal</b> , suggesting that a 50-query budget is	892
		893
		894
		895
		896
		897
		898
		899

Dataset	Budget	ProTeGi	MC	MAPGD
Jailbreak	50	0.833	0.773	0.880
	60	0.836	0.773	0.885
	70	0.842	0.773	0.887
LIAR	50	0.660	0.630	0.710
	60	0.672	0.630	0.720
	70	0.681	0.630	0.725
Sarcasm	50	0.876	0.853	0.913
	60	0.879	0.853	0.918
	70	0.880	0.853	0.920

Table 6: Performance under extended evaluation budgets (60/70 queries). Gains beyond 50 queries are small and do not change the ranking among methods.

Table 7: Test accuracy of MAPGD under different CAAW weighting parameter  $\lambda$ . Performance differences are minor, demonstrating robustness to  $\lambda$ .

$\lambda$	LIAR Accuracy	GSM8K Accuracy
0.5	0.709	0.927
1	0.717	0.935
2	0.711	0.933

sufficient to capture MAPGD’s performance while preserving fairness with prior work (e.g., ProTeGi).

Across these datasets, MAPGD improves by only **+0.007** (Jailbreak), **+0.015** (LIAR), and **+0.007** (Sarcasm) from 50 to 70 queries, confirming a clear diminishing-return trend after approximately 30–40 evaluations, while leaving the main conclusions unchanged.

## D CAAW Parameter Sensitivity

The CAAW module uses a temperature parameter  $\lambda$  to control the emphasis on historically effective agents during gradient fusion. To evaluate the robustness of MAPGD with respect to  $\lambda$ , we conducted experiments on two representative datasets: **LIAR**(binary) and **GSM8K** (mathematical reasoning).

We tested three values of  $\lambda$ : 0.5, 1, and 2. Table 7 summarizes the test accuracy for each setting. The results indicate that moderate variations of  $\lambda$  do not significantly affect the overall performance, supporting the choice of  $\lambda = 1$  in the main experiments.

As shown in Table 7, the test accuracy exhibits only minor variations across the three  $\lambda$  values, with a maximum difference of 0.008 on LIAR and 0.008 on GSM8K. These variations are small relative to the overall performance gains achieved by multi-agent gradient fusion, indicating that while

$\lambda$  influences the sharpness of agent weighting, MAPGD’s performance remains stable for reasonable choices around the default setting of  $\lambda = 1$ .

## E Comparison with Recent Concurrent Work

We are aware of two recent concurrent works that are closely related in spirit to prompt optimization with multi-agent components, namely MAPRO (Zhang et al., 2025) and MA-SAPO (Seo et al., 2025). While these works share a high-level motivation with MAPGD, we did not include them as baselines in our main experimental comparisons due to comparability constraints in problem setting and resource assumptions.

**MAPRO (Zhang et al., 2025).** MAPRO studies prompt optimization *within an explicit multi-agent system (MAS)*: it jointly configures role prompts for multiple agents and leverages the MAS interaction topology (e.g., DAG-structured communication) together with node-/edge-level reward modeling and global inference to select prompt combinations. In contrast, MAPGD targets *task-level prompt optimization* under a fixed query budget, where specialized agents propose pseudo-gradients that are subsequently coordinated (via HCGC and CAAW) to iteratively refine a *single* task prompt. Because MAPRO’s formulation depends on an explicit MAS graph, per-agent prompts, and interaction-specific rewards, directly porting it to our benchmark protocol would require redesigning the environment and supervision signals, making a head-to-head comparison not directly aligned with our evaluation setting.

**MA-SAPO (Seo et al., 2025).** MA-SAPO optimizes prompts through *retrieved reasoning assets*: it constructs a searchable library of semi-structured “reasoning cards” from labeled or scored prompt–response data, and then performs analysis–refinement guided by top- $k$  retrieved assets at test time. This introduces an additional external memory/resource assumption (i.e., an offline asset corpus and its construction pipeline) that is orthogonal to our training-free, black-box optimization protocol. In our evaluation, we intentionally restrict methods to operate under the same query budget without relying on extra annotated corpora or a separately curated retrieval memory. Under these constraints, implementing MA-SAPO would either (i) require additional data collection/annotation to build as-

sets, or (ii) reuse task-specific optimization traces as assets, which changes the method definition and complicates fairness.

**Summary and future work.** Overall, MAPRO and MA-SAPO represent complementary directions—MAS-level joint role-prompt configuration and retrieval-augmented prompt rewriting—whereas MAPGD focuses on multi-agent pseudo-gradient coordination for budgeted task-level prompt refinement. Integrating these paradigms into a unified and strictly comparable evaluation (e.g., extending MAPGD with optional retrieval memories or evaluating MAPRO-style joint prompting on explicit MAS benchmarks) is an important direction for future work.

## F Complexity and Parallelism

The computational bottleneck of MAPGD with HCGC and CAAW (see Algorithm 1) lies in LLM calls for gradient generation, cluster fusion, and prompt expansion, whereas embedding, clustering, similarity checks, and adaptive weighting involve comparatively lightweight vector operations. Multi-agent parallelization significantly reduces wall time, making the framework scalable for practical deployment. In the following, we analyze the per-iteration computational and memory costs, providing a detailed breakdown of the dominant operations and their respective complexities.

**Notation.** Let  $N$  be the number of agents,  $m$  the number of reasons per agent, giving  $G = Nm$  atomic gradients. Embedding dimension  $d$ , clusters  $K \leq G$ , fused gradients  $|\tilde{\mathcal{G}}|$ , expansion variants per gradient  $s$ , MC paraphrases per variant  $n_{\text{mc}}$ , candidate prompts  $|\mathcal{C}|$ , beam width  $k$ , bandit rounds  $T_b$  with  $K_{\text{eval}}$  arms and dev mini-batch size  $b$ .

### Stage costs (time).

1. **Agent gradient generation:** Each of the  $N$  agents generates  $m$  gradients using LLMs, giving  $O(N)$  LLM calls. Parallelization reduces wall time to  $\approx \max t_{\text{LLM}}$ .
2. **HCGC embedding + conflict detection:** Embedding each gradient costs  $O(Gd)$ ; computing pairwise angular conflicts costs  $O(G^2d)$ . For small  $G$  (e.g.,  $G \sim 16$ ), this is negligible compared with LLM calls.
3. **HCGC clustering:** KMeans clustering over  $G$  gradients into  $K$  clusters with  $I$  iterations costs  $O(GKI d)$ , again minor due to small  $G$ .

---

### Algorithm 2 Multi-Agent Textual Gradient Generation

---

**Require:** Agents  $\{A_i\}$ , mini-batch  $M_t$ , task  $T$ , predictor  $\Pi$ , per-agent error cap  $e$ , feedback count  $m$

- 1: **for all** agent  $A_i$  in parallel **do**
  - 2:    $(\hat{y}, y)$  pairs  $\leftarrow T.\text{INFERANDLABEL}(A_i.p, M_t, \Pi)$
  - 3:    $E_i \leftarrow \text{SELECTERRORS}(\hat{y}, y, e)$
  - 4:   **if**  $|E_i| = 0$  **then**  $E_i \leftarrow \text{DIVERSESAMPLES}(M_t, e)$
  - 5:   **end if**
  - 6:    $raw_i \leftarrow \text{LLMGRADIENTPROMPT}(A_i.role, A_i.p, E_i, m)$
  - 7:    $g_i \leftarrow \text{PARSEGRADIENTBLOCKS}(raw_i)$   
    ▷ Split by delimiters
  - 8: **end for**
  - 9: **return**  $\mathcal{G}_t = \{(A_i.role, g_i)\}_{i=1}^N$
- 

4. **HCGC fusion (LLM):** Clusters with multiple gradients require one LLM call per cluster to fuse into a coherent gradient. Total LLM calls up to  $O(K_{\text{merge}})$ , parallelizable across clusters. 1025  
1026  
1027  
1028  
1029
5. **CAAW weighting and fusion:** Per-cluster adaptive weighting is  $O(G)$  vector operations, negligible compared to LLM fusion. 1030  
1031  
1032
6. **Prompt expansion + MC paraphrasing:** Applying each fused gradient:  $O(|\tilde{\mathcal{G}}|)$  gradient applications. MC paraphrasing for each variant:  $O(|\tilde{\mathcal{G}}|sn_{\text{mc}})$  LLM calls. 1033  
1034  
1035  
1036
7. **Diversity filtering:** Computing embeddings:  $O(|\mathcal{C}|d)$ ; naive pairwise similarities:  $O(|\mathcal{C}|^2)$ . For tens of candidate prompts, this remains negligible. 1037  
1038  
1039  
1040
8. **Bandit evaluation:** Probing  $K_{\text{eval}}$  arms over mini-batches of size  $b$  for  $T_b$  rounds:  $O(K_{\text{eval}}bT_b)$ , significantly cheaper than exhaustive evaluation  $O(|\mathcal{C}||D_{\text{dev}}|)$ . 1041  
1042  
1043  
1044

**Space complexity.** Text storage:  $O(|\mathcal{C}|L_{\text{avg}})$ . Embeddings:  $O((G + |\mathcal{C}|)d)$ , typically a few MB. Optional caches (unique prompts or intermediate LLM outputs) scale linearly with the number of prompts and clusters. 1045  
1046  
1047  
1048  
1049

## G Theoretical Analysis

In this section, we provide convergence guarantees for MAPGD under mild assumptions, following the 1050  
1051  
1052

---

**Algorithm 3** Hypersphere-Constrained Gradient Clustering
 

---

**Require:** Agent gradients  $\mathcal{G}_t$ 

- 1: Embed and normalize:  $\hat{v}_k = \phi(g_k) / \|\phi(g_k)\|$
  - 2: Detect conflicts:  $\mathcal{C}_{\text{conf}} \leftarrow \{(i, j) : \cos^{-1}(\hat{v}_i^\top \hat{v}_j) > \theta\}$
  - 3: Cluster gradients:  $\{S_k\}_{k=1}^K \leftarrow \text{KMEANS}(\{\hat{v}_k\}, K_{\text{max}})$
  - 4: **for**  $k = 1$  to  $K$  **do**
  - 5:   Apply angular margin: reassign gradients violating  $\cos(n \cdot \Delta(\hat{v}_i, c_k)) > \cos(\Delta(\hat{v}_i, c_j))$
  - 6:   Fuse cluster  $S_k$  via LLM:  $f_k \leftarrow \text{LLMFUSE}(S_k, \mathcal{C}_{\text{conf}})$
  - 7: **end for**
  - 8: **return** Fused clusters:  $\tilde{\mathcal{G}}_t = \{f_1, \dots, f_K\}$
- 

**Algorithm 4** Channel-Adaptive Agent Weighting and Fusion
 

---

**Require:** Clustered gradients  $\tilde{\mathcal{G}}_t$ , dev data  $D_{\text{dev}}$ , temperature  $\lambda$ 

- 1: **for all** cluster  $\tilde{\mathcal{G}}_t$  **do**
  - 2:   Compute per-gradient validation gain  $s_i$  on  $D_{\text{dev}}$
  - 3:   Assign adaptive weight:  $w_i = \frac{\exp(\lambda s_i)}{\sum_{j \in S_k} \exp(\lambda s_j)}$
  - 4:   Fuse cluster gradients:  $f_k = \Psi\left(\sum_{i \in S_k} w_i g_i\right)$
  - 5: **end for**
  - 6: **return**  $F_t = \{f_1, \dots, f_K\}$
- 

1053 stochastic approximation framework. Our goal is  
 1054 to bridge the gap between the continuous optimization  
 1055 theory of stochastic gradient descent (SGD)  
 1056 and the discrete prompt optimization carried out  
 1057 in MAPGD. We show that, despite operating in  
 1058 a structured and discrete search space, MAPGD  
 1059 achieves the same sublinear convergence rate of  
 1060  $O(1/\sqrt{T})$  in both convex and non-convex settings.

### 1061 G.1 Assumptions

1062 We begin with a set of assumptions standard in  
 1063 stochastic optimization but reinterpreted in the con-  
 1064 text of multi-agent prompt optimization.

- **(A1) Alignment (Unbiasedness).** For some  $\mu > 0$ , the stochastic semantic gradient  $g^{(t)}$  maintains alignment with the true gradient:

$$1065 \mathbb{E}\left[\langle g^{(t)}, \nabla F(p^{(t)}) \rangle \mid p^{(t)}\right] \geq \mu \|\nabla F(p^{(t)})\|^2.$$

---

**Algorithm 5** SynchronizeAgents
 

---

**Require:** Agents  $\{A_i\}_{i=1}^N$ , current best prompt  $p_\star^{(t)}$ 

- 1: **for** each agent  $A_i$  **do**
  - 2:    $A_i.p \leftarrow p_\star^{(t)}$    ▷ Set the agent’s current prompt to the global best
  - 3:   RESETGRADIENTHISTORY( $A_i$ )   ▷ Optional: clear outdated gradient history
  - 4:   UPDATEPERFORMANCEMEMORY( $A_i, p_\star^{(t)}$ ) ▷ Record performance feedback for the new prompt
  - 5: **end for**
  - 6: **return** Updated agents  $\{A_i\}$
- 

This reflects the role of semantic fusion: multi-agent aggregation reduces the chance of adversarial or noisy updates, ensuring progress along descent directions.

- **(A2) Bounded Second Moment.** For constants  $\rho, \sigma^2 \geq 0$ ,

$$\mathbb{E}\left[\|g^{(t)}\|^2 \mid p^{(t)}\right] \leq \rho \|\nabla F(p^{(t)})\|^2 + \sigma^2.$$

This captures the variance-control effect of the bandit-based selection mechanism, which prevents uncontrolled explosion of gradient magnitude.

- **(A3) Smoothness or Lipschitzness.** For convex tasks,  $F$  is  $G$ -Lipschitz with domain diameter  $D$ . For non-convex tasks,  $F$  is  $L$ -smooth:  $\|\nabla F(u) - \nabla F(v)\| \leq L\|u - v\|$ .

### G.2 Supporting Lemmas

We restate two standard lemmas, adapted to the MAPGD setting.

**Lemma 1 (Convex Projection Inequality).** For convex  $F$  with feasible set  $\mathcal{P}$ , the projected subgradient update

$$p^{(t+1)} = \Pi_{\mathcal{P}}\left(p^{(t)} - \eta g^{(t)}\right)$$

satisfies

$$\|p^{(t+1)} - p^*\|^2 \leq \|p^{(t)} - p^*\|^2 - 2\eta \langle g^{(t)}, p^{(t)} - p^* \rangle + \eta^2 \|g^{(t)}\|^2.$$

**Lemma 2 (Non-Convex Descent Lemma).** If  $F$  is  $L$ -smooth, then for update  $p^{(t+1)} = p^{(t)} - \eta g^{(t)}$ , we have

$$F(p^{(t+1)}) \leq F(p^{(t)}) - \eta \langle \nabla F(p^{(t)}), g^{(t)} \rangle + \frac{L}{2} \eta^2 \|g^{(t)}\|^2.$$

### 1097 G.3 Main Results

1098 **Convex Convergence.** Suppose  $F$  is convex,  $G$ -  
 1099 Lipschitz, and  $\mathcal{P}$  has diameter  $D$ . Let  $\bar{p}_T =$   
 1100  $\frac{1}{T} \sum_{t=1}^T p^{(t)}$ . Under (A1)–(A2) and step size  $\eta =$   
 1101  $\frac{D}{G\sqrt{T}}$ , we obtain:

$$1102 \mathbb{E}[F(\bar{p}_T)] - F(p^*) = O\left(\frac{1}{\sqrt{T}}\right).$$

1103 *Proof sketch.* By Lemma 1 and convexity:

$$1104 F(p^{(t)}) - F(p^*) \leq \langle g^{(t)}, p^{(t)} - p^* \rangle.$$

1105 Summing over  $t = 1, \dots, T$  and applying (A1)–  
 1106 (A2), we bound the regret:

$$1107 \sum_{t=1}^T \mathbb{E}[F(p^{(t)}) - F(p^*)] \leq \frac{D^2}{2\eta} + \frac{\eta G^2 T}{2}.$$

1108 Using Jensen’s inequality for  $\bar{p}_T$  and optimizing  $\eta$ ,  
 1109 we conclude the  $O(1/\sqrt{T})$  rate.

1110 **Non-Convex Convergence.** Suppose  $F$  is  $L$ -  
 1111 smooth and (A1)–(A2) hold. With constant step  
 1112 size  $\eta = \Theta(1/\sqrt{T})$ , we have

$$1113 \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\nabla F(p^{(t)})\|^2] = O\left(\frac{1}{\sqrt{T}}\right).$$

1114 *Proof sketch.* Applying Lemma 2 and taking  
 1115 conditional expectation:

$$1116 \mathbb{E}[F(p^{(t+1)})] \leq \mathbb{E}[F(p^{(t)})] - \eta\mu\mathbb{E}[\|\nabla F(p^{(t)})\|^2]  
 1117 + \frac{L}{2}\eta^2 \left( \rho\mathbb{E}[\|\nabla F(p^{(t)})\|^2] + \sigma^2 \right). \quad (12)$$

1118 Summing over  $t = 1 \dots T$  gives

$$1119 \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\nabla F(p^{(t)})\|^2] \leq \frac{2(F(p^{(1)}) - F_{\text{inf}})}{\mu T \eta} \\ + \frac{L\sigma^2}{\mu}\eta.$$

1120 Balancing terms with  $\eta = \Theta(1/\sqrt{T})$  yields the  
 1121 claimed rate.

### 1122 G.4 Summary of Theoretical Guarantees

1123 In this work, we bridge the gap between classical  
 1124 stochastic optimization, which assumes continu-  
 1125 ous parameter spaces, and the inherently discrete,  
 1126 structured nature of prompt optimization. Our con-  
 1127 vergence analysis for MAPGD does not require the  
 1128 discrete prompt space to be continuous. Instead,

1129 we show that the framework’s core mechanisms  
 1130 ensure that the optimization process satisfies the  
 1131 key conditions of stochastic approximation theory.

The argument unfolds along three main dimen-  
 1132 sions: 1133

- **Directional Alignment (A1):** Although tex-  
 1134 tual pseudo-gradients are not true mathemati-  
 1135 cal gradients, HCGC aggregates semantically  
 1136 coherent suggestions from multiple agents  
 1137 into a fused gradient that statistically aligns  
 1138 with the true descent direction. 1139

- **Variance Control (A2):** Discrete prompt ed-  
 1140 its are prone to high variance, but the bandit-  
 1141 based candidate selection mechanism filters  
 1142 out unreliable or detrimental candidates. This  
 1143 effectively bounds the second moment of  
 1144 stochastic updates. 1145

- **Smoothness Justification (A3):** We assume  
 1146 the empirical loss function is Lipschitz or  
 1147 smooth in the semantic embedding space,  
 1148 even though the prompts themselves are dis-  
 1149 crete. 1150

Given that these conditions are met, MAPGD  
 1151 inherits the convergence guarantees of classical  
 1152 stochastic gradient descent. Consequently, despite  
 1153 operating in a discrete textual space, MAPGD  
 1154 achieves a robust sublinear convergence rate of  
 1155  $O(1/\sqrt{T})$ , providing a solid theoretical founda-  
 1156 tion for the observed empirical effectiveness of our  
 1157 multi-agent, geometry-aware approach to prompt  
 1158 optimization. 1159

## 1160 H Qualitative Case Study: From 1161 Conflicting Pseudo-Gradients to 1162 Prompt Updates

### 1163 H.1 Case Setup

1164 We select a representative optimization step on  
 1165 the LIAR task where multiple specialized agents  
 1166 propose semantically conflicting pseudo-gradients.  
 1167 This step exhibits clear disagreement among agents  
 1168 while still leading to a measurable improvement  
 1169 in downstream validation behavior after gradient  
 1170 fusion.

Item	Value
Task	LIAR
Iteration	$t = 3$
Agents	4
Raw Gradients	16
Illustrated Conflict Pairs	2
HCGC Clusters	3

At this iteration, agents disagree on whether errors primarily stem from ambiguous label semantics, missing contextual cues, or insufficient decision boundary clarification, making it a suitable case for qualitative inspection.

## H.2 Raw Pseudo-Gradients from Specialized Agents

We focus on two representative conflict pairs extracted from the same optimization step, which illustrate structurally different failure diagnoses and motivate the need for explicit conflict handling.

### Conflict Pair 1: Format Clarity vs Example Diversity. Format Designer:

<START> The prompt instructs the model to answer with ‘Yes’ or ‘No’, but it does not specify what these labels represent. This ambiguity can cause the model to misinterpret the classification objective. <END>

#### Example Curator:

<START> The provided examples do not cover a sufficiently diverse range of statements. Including more varied and representative cases could help the model generalize better. <END>

The cosine similarity between these two gradients is 0.13 (angle = 1.44 rad), indicating strong semantic disagreement. Importantly, the conflict is not superficial: the format-oriented gradient advocates tightening the input–output contract through explicit label semantics, while the example-oriented gradient pushes the prompt toward structural expansion via additional demonstrations. Applying both updates independently would lead to incompatible prompt modifications—one emphasizing deterministic label interpretation and the other reallocating prompt capacity to example coverage.

### Conflict Pair 2: Example Coverage vs Definition of “Lie”. Example Curator:

<START> The examples fail to illustrate borderline or ambiguous cases, which may limit the model’s ability to generalize across different claim types. <END>

### Instruction Specialist:

<START> The prompt does not clearly define what constitutes a ‘lie’ versus an inaccurate or exaggerated statement, leading to inconsistent decisions. <END>

This pair exhibits a cosine similarity of 0.23 (angle = 1.34 rad), reflecting a fundamental disagreement about the source of errors. The example-focused gradient assumes that misclassification primarily stems from insufficient empirical coverage, whereas the instruction-focused gradient attributes failures to an underspecified decision criterion. These two perspectives induce different optimization directions: expanding illustrative evidence versus formalizing the semantic boundary of the task itself.

**Discussion.** Both conflict pairs arise from valid but incompatible diagnostic assumptions. Crucially, neither gradient is inherently incorrect; each targets a distinct aspect of model behavior. This highlights the core challenge addressed by MAPGD: without explicit conflict detection and coordination, naïvely aggregating such gradients would obscure their semantic tension and lead to unstable or diluted prompt updates.

### H.3 Conflict Detection and HCGC Clustering

Based on the two conflict pairs in §H.2, we observe clear semantic disagreement among agent-proposed directions. In particular, the *Format Clarity vs Example Diversity* pair has cosine similarity 0.13 (angle = 1.44 rad), and the *Example Coverage vs Definition of “Lie”* pair has cosine similarity 0.23 (angle = 1.34 rad). The corresponding angles (1.44 and 1.34 rad) indicate substantial directional mismatch, and both similarities fall below the conflict threshold. HCGC therefore marks them as conflicting and avoids naïve averaging.

After projecting all gradient embeddings onto the unit hypersphere and applying the angular-margin constraints, HCGC partitions the 16 raw pseudo-gradients into three coherent clusters, as summarized in Table 8.

This clustering reflects a key property of HCGC: conflicting gradients are not forced into a single averaged direction. Instead, they are separated into semantically consistent groups (Table 8), enabling subsequent fusion to preserve complementary information (e.g., clarifying label semantics in C1 while retaining coverage-oriented refinements in C2) without letting one channel dominate the other.

Table 8: HCGC clustering summary for the representative LIAR optimization step ( $t=3$ ).

Cluster	Agents Involved	Semantic Theme
C1	instruction, format	label semantics and task definition (binding Yes/No to True/False, clarifying what is being judged)
C2	example, style	coverage and contextual cues for borderline cases (exaggeration, unverifiable framing, pragmatic signals)
C3	instruction only	decision-boundary calibration (strictness vs recall; discouraging overly permissive acceptance of absolute claims)

The semantic themes are manually summarized by the authors to reflect the dominant intent of each cluster, rather than being directly emitted by the model.

#### H.4 Channel-Adaptive Agent Weighting

Given the clustered directions above, CAAW assigns each agent a contribution weight based on its recent performance signal, so that the subsequent fusion emphasizes gradients that are empirically reliable for the current task state. At this iteration, the adaptive weights are:

$$w_{\text{instruction}} = 0.41, \quad w_{\text{example}} = 0.33, \\ w_{\text{format}} = 0.17, \quad w_{\text{style}} = 0.09.$$

These weights are consistent with the observed conflict structure. Clusters C1 and C2 correspond to two *structurally competing* but both valid refinement axes: (i) clarifying label semantics and the task definition (C1), and (ii) improving coverage and contextual handling of borderline cases (C2). CAAW prioritizes instruction- and example-driven gradients while down-weighting stylistic adjustments, preventing prompt capacity from being consumed by surface-level phrasing when the dominant failures are semantic (label meaning) and epistemic (what constitutes a lie).

#### H.5 Fused Gradient and Prompt Modification

After separating conflicting directions with HCGC and weighting them with CAAW, MAPGD fuses the cluster-level signals into a single coherent update that preserves complementary intent across clusters. In this case, the fused direction primarily integrates (a) C1: explicit label semantics and task definition, and (b) C2: guidance for borderline or unverifiable claims, while avoiding excessive prompt expansion.

**Before:**

“The following statement is either true or false. Answer Yes or No.”

**After:**

“Determine whether the statement is *factually accurate*. Treat Yes as True and No as False. If the statement relies on exaggeration, satire, or unverifiable claims, judge it as False. Return only Yes or No.”

This modification reflects the fused gradient in two concrete ways: (i) it resolves the *format/label ambiguity* highlighted in C1 by explicitly binding Yes/No to True/False, and (ii) it incorporates the *borderline-case* handling emphasized by C2 without requiring lengthy demonstrations, thereby avoiding the prompt-structure conflict that would arise from naively mixing rule-heavy and example-heavy updates.

#### H.6 Behavioral Change on Previously Misclassified Examples

We report qualitative behavioral changes on examples that were previously misclassified due to (i) label ambiguity and (ii) overly permissive interpretation of absolute or exaggerated claims.

**Example 1 (Absolute claim / overgeneralization).**

**Claim:** “The senator has voted against every environmental bill.”

**Before:** Yes      **After:** No

The updated prompt instructs the model to treat unverifiable or sweeping claims as False, which directly corrects the previous tendency to accept absolute statements (e.g., “every”) as True.

**Example 2 (Borderline / unverifiable framing).**

**Claim:** “The policy has been a complete disaster for all families.”

**Before:** Yes      **After:** No

This change follows the newly added rule that exaggeration or unverifiable claims should be judged as False, reducing false positives driven by rhetorical or hyperbolic language (e.g., “complete disaster”).

## I Case Study: System Prompt Optimization

To further illustrate the applicability of MAPGD, we present a case study where our method is applied to optimize the system prompt of a large language model assistant. The original prompt is designed to support multi-source analysis and financial data interpretation for decision-making tasks. Using MAPGD, we refine the prompt to enhance robustness, accuracy, and interpretability by embedding explicit verification protocols, structured analysis guidelines, and risk prioritization frameworks. This example highlights how MAPGD can be deployed in practical LLM applications beyond benchmark datasets, particularly in domains where data authenticity, reliability, and interpretability are critical.

**Original System Prompt.** You are an AI assistant designed to process, analyze, and synthesize information from multiple sources in order to answer user questions, generate insights, and prepare detailed reports. You have specialized capabilities in financial data interpretation, knowledge retrieval, and multi-source analysis. You support both operational and strategic decision-making for Golden Section’s portfolio companies.

**Optimized System Prompt via MAPGD.** You are an AI assistant specialized in processing, analyzing, and synthesizing information from multiple sources to answer user questions, generate insights, and prepare detailed reports. Your core capabilities include financial data interpretation, knowledge retrieval, and multi-source analysis, with a focus on supporting operational and strategic decision-making for Golden Section’s portfolio companies.

In performing your duties, you must ensure that all input data undergoes rigorous verification for authenticity, accuracy, and completeness before any analysis is conducted. This includes implementing protocols to validate financial figures, legal terms, and other critical information for alignment with established norms and credible sources. Your process must involve:

- Cross-referencing information from diverse, credible sources to detect and mitigate false, exaggerated, or incomplete data.
- Assessing the reliability of each source, prioritizing primary sources where available.
- Identifying and resolving inconsistencies, ambiguities, or potential misinformation through

systematic checks.

- Ensuring all risk assessments and conclusions are based solely on validated and accurate inputs to maintain the integrity of your outputs.

### Data Authenticity and Completeness Verification

- Scrutinize contextual cues (e.g., “Context: Section: Payback Period:”) to ensure alignment with expected data types and structures.
- Check for numerical or factual inconsistencies, such as typos (e.g., “11975” instead of “1975”), exaggerations (e.g., “\$22.0M” without supporting context), or missing critical information.
- Validate that all referenced data points are present, logically consistent, and contextually appropriate.
- Flag and document any anomalies for further investigation before proceeding with classification or analysis.

### Context Interpretation and Parsing Guidelines

- Carefully interpret and utilize contextual cues, especially in nested or ambiguous contexts (e.g., “Context: Section: Name & Headquarters:”).
- Accurately parse section headers and contextual clues to prevent misclassification or incomplete analysis in multi-section reports.
- Anchor analysis to the document’s structure by adhering to hierarchical or sequential organization.
- Resolve discrepancies in contextual labeling or structure to maintain coherence.

### Structured Classification and Risk Prioritization Framework

1. **Factual Reporting and Descriptive Analysis:** Present verified information such as corporate history, operational metrics, and financial data neutrally, before transitioning to evaluative content.

1430	2. <b>Business Analysis:</b> Evaluate performance, market positioning, and strategic initiatives; assess risks by severity, likelihood, and propose contextualized mitigation.	<b>example_curator</b> Focuses on selecting representative, diverse examples and ensuring consistent formatting.	1476
1431			1477
1432			1478
1433			
1434	3. <b>Legal Risk Analysis:</b> Examine compliance, regulatory, and contractual risks; assess impact and propose mitigation actions aligned with legal context.	<b>format_designer</b> Designs clear output templates and structured formats for better model understanding.	1479
1435			1480
1436			1481
1437		<b>style_optimizer</b> Optimizes language expression for professionalism and task-specific adaptation.	1482
1438			1483
1439	4. <b>Cross-Domain Analysis:</b> For overlapping elements, classify by primary context and document dual-category cases with rationale.		1484
1440			
1441	<b>Validation Mechanisms for Cross-References</b>	<b>J.2 Agent Roles for Mathematical Reasoning Tasks</b>	1485
1442	• Distinguish between source types (e.g., governance vs. identity records).	For mathematical reasoning, the agent roles are adapted to focus on logical decomposition, calculation, and problem interpretation.	1486
1443			1487
1444	• For each statement, explicitly identify the source type and ensure contextual alignment.	<b>reasoning_specialist</b> Specializes in enhancing mathematical reasoning processes, ensuring clear step-by-step logic and proper problem decomposition.	1488
1445			1489
1446	<b>Empirical Observation.</b> After applying MAPGD, we observe a notable increase in the comprehensiveness of the generated financial analysis reports. Specifically, the average report length increased from 1616 words under the original system prompt to 1965 words with the optimized prompt. This indicates that the optimized prompt encourages the model to produce more detailed and contextually grounded content. While the per-report output length becomes longer, the improved accuracy and completeness reduce the need for repeated generations or manual corrections, thereby lowering the overall token consumption in practical workflows. These findings further support that MAPGD enhances the robustness and effectiveness of system prompts in real-world generation tasks.	<b>calculation_optimizer</b> Focuses on improving calculation accuracy, suggesting better computational methods, and ensuring numerical correctness.	1490
1447			1491
1448			1492
1449		<b>problem_interpreter</b> Specializes in interpreting math word problems, extracting relevant information, and identifying mathematical relationships.	1493
1450			1494
1451			1495
1452			1496
1453			1497
1454		<b>solution_formatter</b> Optimizes mathematical solution presentation, ensures clear formatting and proper answer notation (e.g., #### [final_answer]).	1498
1455			1499
1456			1500
1457			1501
1458		<b>J.3 Agent Roles for Financial Report Analysis</b>	1502
1459		For the complex task of business and financial analysis, agents are given highly specialized roles focusing on data verification, structured report generation, and professional tone.	1503
1460			1504
1461			1505
1462	<b>J Agent Specialization for Different Tasks</b>	<b>instruction_specialist</b> Crafts clear, executable instructions for generating truthful, structured reports and rejecting false or exaggerated inputs. Instructions must include steps to verify input data authenticity (e.g., checking financial metrics for realism).	1511
1463	To adapt MAPGD to different domains, we define specialized agent roles tailored to the unique requirements of each task. Below are the agent configurations used for the classification, mathematical reasoning, and financial analysis tasks in our experiments.		1512
1464			1513
1465			1514
1466			1515
1467			1516
1468			
1469	<b>J.1 Agent Roles for Classification Tasks</b>	<b>example_curator</b> Selects and formats diverse examples that demonstrate truthful report generation and false data rejection, including both positive (realistic) and negative (exaggerated) samples.	1517
1470	For general classification tasks, the agents focus on the core components of a good prompt: instructions, examples, format, and style.		1518
1471			1519
1472			1520
1473	<b>instruction_specialist</b> Analyzes and improves task instructions, ensuring clarity, completeness, and executability.		1521
1474			
1475			

1522 **format\_designer** Specializes in designing clear,  
1523 structured output templates for reports (e.g.,  
1524 using <Business Analysis>, <Financial  
1525 Report> tags) and error messages (e.g.,  
1526 <Error>).

1527 **style\_optimizer** Optimizes language style and  
1528 tone for generating professional reports in  
1529 business, financial, and legal contexts, while  
1530 ensuring error messages are direct and unam-  
1531 biguous.

## 1532 K Example Prompts Before and After 1533 Optimization

1534 This appendix provides the initial prompts and the  
1535 final optimized prompts generated by MAPGD for  
1536 each of the benchmark datasets used in our experi-  
1537 ments.

### 1538 K.1 LIAR Dataset

#### Initial Prompt

```
# Task
Determine whether the Statement is a lie
(Yes) or not (No) based on the Context
and other information.
# Output format
Answer Yes or No as labels
# Prediction
Text: {text}
label:
breakindent
```

#### Optimized Prompt

```
# Task
Evaluate the given Statement using the
provided context.

Classify as 'Yes' if the Statement is
intentionally false
or misleading (i.e., a lie).

Classify as 'No' if the Statement is true,
or if it is
incorrect due to a mistake or lack of
knowledge rather
than intentional deception.

# Output format
Output exactly one label: 'Yes' or 'No'.
'Yes' indicates the Statement is a lie; 'No'
indicates
it is not a lie.

# Prediction
Text: {text}
label:
breakindent
```

### K.2 Ethos Dataset

#### Initial Prompt

```
# Task
Is the following text hate speech?
# Output format
Answer Yes or No as labels
# Prediction
Text: {text}
Label:
breakindent
```

#### Optimized Prompt

```
# Task
Your task is to analyze the following text
and determine if
it contains hate speech. For the purpose of
this task, hate
speech is defined as any speech, gesture,
conduct, writing,
or display that may incite violence or
prejudicial actions
against or by a protected individual or
group, or because it
disparages or intimidates a protected
individual or group.
Please base your judgment on the explicit
content of the text,
rather than personal interpretations or
assumptions.
# Output format
Provide your answer in a clear format as
either
'Contains Hate Speech' or 'Does Not Contain
Hate Speech'
# Prediction
Text: {text}
Label:
breakindent
```

### K.3 Jailbreak Dataset

#### Initial Prompt

```
# Task
Detect if the message is a jailbreak attack,
i.e. an attempt by a user to break through an
AI system's protections
# Output format
Answer Yes or No as labels
# Prediction
Text: {text}
Label:
breakindent
```

#### Optimized Prompt

```
# Task
Classify the message as one of the following
labels based on its potential to be a
jailbreak attack, i.e., an attempt to
bypass an AI system's protections:
- "OvertAttack": The message contains clear,
direct attempts
to break through system safeguards.
- "SuspiciousContext": The message shows
indirect or
```

1540

contextual cues that may indicate a multi-stage or sophisticated attack attempt.

- "Benign": The message poses no apparent risk and does not attempt to circumvent protections.

# Output format  
Select and output exactly one label from the options above.

# Prediction  
Text: {text}  
Label:  
breakindent

4. Select the correct answer (A, B, C, D, or E)

# Output Format  
Show your reasoning and end with: Answer: [LETTER]  
(For example: "Answer: A" or "Answer: B")

# Problem  
{text}

# Options  
{options}  
breakindent

#### K.4 GSM8K Dataset

##### Initial Prompt

# Task  
Solve the math word problem step by step.  
# Instructions  
1. Read the problem carefully  
2. Identify what needs to be calculated  
3. Show your work step by step  
4. Provide the final numerical answer  
# Output Format  
Show your reasoning process and end with: #### [final\_answer]  
# Problem  
{text}  
breakindent

##### Optimized Prompt

# Task  
Solve the math word problem step by step.  
# Instructions  
1. Read the problem carefully  
2. Identify what needs to be calculated  
3. If there are ambiguous situations in the problem, make reasonable assumptions and state them clearly  
4. Show your work step by step  
5. Provide the final numerical answer  
# Output Format  
Show your reasoning process, including any assumptions made, and end with: #### [final\_answer]  
# Problem  
{text}  
breakindent

#### K.5 AQUARAT Dataset

##### Initial Prompt

# Task  
Solve the math word problem and choose the correct answer from the given options.

# Instructions  
1. Read the problem carefully  
2. Analyze each option  
3. Show your reasoning step by step

##### Optimized Prompt

# Task  
Solve the provided problem by systematically classifying its type, analyzing all constraints, applying formal mathematical reasoning, and verifying the solution against all conditions before selecting the correct answer.

# Instructions  
Adhere strictly to the following structured framework:

- \*\*Problem Classification:\*\***
  - \* Explicitly identify the core problem domain (e.g., "Combinatorial Selection," "Constraint Satisfaction," "Partnership Profit-Sharing," "Work Rate").
  - \* Justify the classification with a brief rationale based on the problem statement.
- \*\*Constraint Analysis:\*\***
  - \* Enumerate all explicit constraints from the problem text.
  - \* Infer and list any implicit constraints or logical dependencies.
  - \* Categorize each constraint (e.g., Hard/Must-Fulfill, Soft/Optimization) and specify their logical relationships (e.g., AND, OR).
- \*\*Mathematical Formalization & Numerical Context Resolution:\*\***
  - \* Identify and resolve any ambiguous numerical references (e.g., "monthly" vs. "annual", "together" vs. "individual").
  - \* Infer and incorporate any implicit numerical variables (e.g., initial quantities, unstated rates, starting points).
  - \* Normalize all units to a consistent basis for calculation.

- \* Define all variables, sets, and parameters using clear notation (e.g.,  $C_{total}$ ,  $P_A$ ,  $S_{eligible}$ ,  $nCr$ ).
- \* Structure the problem using appropriate formalisms: timelines, sets, equations, or logical statements.
- \* For combinatorial problems, explicitly state the formula used.

4. **Solution Execution**

- \* Perform all calculations step-by-step, showing substitutions and operations.
- \* For combinatorial scenarios, enumerate valid combinations or calculate cardinalities.
- \* Derive the target value (e.g., profit share, number of committees, optimal value).

5. **Verification & Solution Finalization**

- \* Systematically verify that the proposed solution satisfies every constraint from Step 2.
- \* For combinatorial answers, confirm counts against constraints using direct checks or complementary counting.
- \* Match the final, verified result to the provided options to select the correct answer.

**Output Format**  
Your final output must follow this structure precisely:

**Problem Classification**  
[Your classification and rationale]  
breakindent

reasoning. Your primary goal is to correctly parse and compute multi-step problems involving distinct object categories and their quantitative relationships.

**Instructions**

1. **Parse and Categorize**  
Read the problem and question carefully. Identify and list all distinct object types or categories (e.g., 'packages of gum', 'boxes of candy'). For each category, extract:
  - The number of units (e.g., 5 packages).
  - The quantity per unit, if specified (e.g., 8 pieces per package).
2. **Structure the Solution Hierarchically**
  - **Step 1: Calculate Intra-Category Totals**  
For each identified category, calculate its total quantity. If a unit rate is given, perform the multiplication (e.g.,  $5 \text{ packages} * 8 \text{ pieces/package} = 40 \text{ pieces of gum}$ ). If no unit rate exists, use the given quantity directly.
  - **Step 2: Perform Inter-Category Operations**  
Using the totals from Step 1, now perform the final operation as required by the question (e.g., sum all category totals for a grand total, or find the difference between two category totals for a comparison).
3. **Show Your Work**  
Present your reasoning clearly, reflecting this two-step hierarchical process. First, show all calculations within categories. Then, show the final combination of these category totals.
4. **Final Answer**  
Box your final numerical result.

**Output Format**  
Reasoning: [Your step-by-step reasoning here]  
#### [final\_answer]

**Problem**  
{text}

**Question**  
{question}  
breakindent

## K.6 SVAMP Dataset

### Initial Prompt

**Task**  
Solve the math word problem step by step.

**Instructions**

1. Read the problem carefully
2. Identify what needs to be calculated
3. Show your work step by step
4. Provide the final numerical answer

**Output Format**  
Show your reasoning process and end with:  
#### [final\_answer]

**Problem**  
{text}  
breakindent

### Optimized Prompt

**Task**  
Solve the math word problem through structured, hierarchical